# NVIDIA Ampere vGPU 配置与测试参考

修订记录

Date	Revision	Editor	Changes
2020/11/12	1.0	Merlin Ma	Initial version.
2020/11/16	1.1	Merrin Ma	Some typo corrections and CI testing.
2021/08/05	2.0	Merlin Ma	For Ampere GPU models (not only A100).
			None-MIG SR-IOV vGPU configuration added. Document
			restructuring. (chapter 4)
2021/08/10	2.1	Merlin Ma	Detailed SR-IOV GPU explanation, rules (chapter 2.4, 2.5).
			More examples.

#### 概述

NVidia A100 GPU 提供了全新的 MIG 多实例 GPU 模式, MIG 功能可以将 GPU 安全地划分 成为最多 7 个独立地 GPU 计算实例,每个计算实例拥有自己所属的计算资源,提供更好 的 QoS 以及故障隔离能力。从 vGPU 11.1 版本开始支持基于 MIG 技术的 vGPU 虚拟机实 例。本文将着重介绍:使用 NVidia A100 GPU 在 MIG 模式下 vGPU 实例的创建和管理,主要 面向 MIG vGPU 的 KVM 虚拟化环境 PoC 测试和部署方法参考。同时,本文也以 A100 GPU 为例,介绍了基于 SR-IOV 的 GPU 的 vGPU 配置方法,也适用于 A10、A40、A16 等 GPU 的 时分 vGPU 方案。

## 1. 安装需求和准备工作

- 1. 主机配置: NVidia Ampere GPU 的 X86\_64 服务器。
- 2. KVM HOST: 安装 Redhat RHEL 8.2 及其以上版本。
- 3. GRID 11.1 或者以上版本试用软件及许可证,请先确认 vGPU 软件与物理 GPU 的兼容性。

## 2. KVM Host 软硬件基础配置概要

#### 2.1. 修改 Linux kernel 启动参数

在 KVM 平台的 Linux kernel 启动项中添加 intel\_iommu=on iommu=pt 通常此配置文件在: /boot/efi/EFI/redhat/grubenv 或 /boot/grub2/grubenv

这是修改后的 RHEL8.2 kenel 参数: (添加蓝色部分):

[root@kvm ~]# cat /boot/grub2/grubenv
# GRUB Environment Block
saved entry=11f957e144e643c68255fd47642d77b3-4.18.0-193.el8.x86 64
kernelopts=root=/dev/mapper/rhel-root ro crashkernel=auto resume=/dev/mapper/rhel-swap rd.lvm.lv=rhel/root rd.lv
m.lv=rhel/swap rhgb intel_iommu=on iommu=pt quiet
boot_success=0
boot_indeterminate=0
[root@kvm ~]# cat /proc/cmdline
BOOT_IMAGE=(hd2,gpt2)/vmlinuz-4.18.0-193.el8.x86_64 root=/dev/mapper/rhel-root ro crashkernel=auto resume=/dev/m
apper/rhel-swap rd.lvm.lv=rhel/root rd.lvm.lv=rhel/swap rhgb intel_iommu=on iommu=pt quiet
[root@kvm ~]#

## 2.2. 安装 vGPU 的 KVM Host 驱动

rpm -ivh NVIDIA-vGPU-rhel-8.2-450.89.x86\_64.rpm。

[noot@kym lloct]# ]c ] NV/TDTA vCDU pho	2 2 4 50 90 90 90 50 51 000	
[POOL@KVM HOSC]# IS -I NVIDIA-VGPO-PME	er-8.2-450.89.886_64.1.hill	
-rw-rw-r 1 root root 18239080 Nov 5	5 01:13 NVIDIA-vGPU-rhel-8.2-450.89	.x86_64.rpm
[root@kvm Host]# rpm -ivh NVIDIA-vGPU-	rhel-8.2-450.89.x86_64.rpm	
Verifying	****	[100%]
Preparing	****	[100%]
Updating / installing		
1:NVIDIA-vGPU-rhel-1:8.2-450.89	****	[100%]
[root@kvm Host]# <mark> </mark>		

## 2.3. Reboot 并确认 GPU 状态

修改 kernel 参数和安装 vGPU Host 驱动以后,一定要 reboot 使其生效。

检查 nvidia-smi 输出: 应该看到 A100 GPU, 默认状态 MIG 是被禁用状态的。 记录下 A100 GPU 的 PCI BUS ID, 本例中为: 0000:81:00.0

[root@ Thu No	kvm ~]# v 12 23:	nvid 17:1	ia-smi 9 2020						
+   NVID	IA-SMI 4	450.8	9	Driver	Version:	450.8	89 (	CUDA Versio	on: N/A
GPU   Fan 	Name Temp F	Perf	Persist Pwr:Usa	ence-M ge/Cap	Bus-Id	Memor	Disp.A ry-Usage	Volatile   GPU-Util 	Uncorr. ECC Compute M. MIG M.
0   N/A 	A100-P0 29C	CIE-4 P0	0GB 35W /	On 250W	0000000 0M:	9:81:0 iB / 4	00.0 Off 40537MiB	     0%	0 Default Disabled
1   N/A 	Tesla 1 35C	Γ4 P8	16W /	On 70W	0000000 83M:	0:84:0 iB / 1	00.0 Off 15359MiB	     0%	0   Default   N/A

#### 2.4. 启用 SRIOV

注意:在执行此步骤之前,请确保物理 GPU 未被任何其他进程使用,例如: CUDA 程序、 监控程序或 nvidia-smi 命令。仅使用随 NVIDIA vGPU 软件提供的自定义脚本 sriov-manage。 不要尝试通过其他方式使能 SR-IOV。

运行/usr/lib/nvidia/sriov-manage -e 加上前面记下的 Bus-ID, 例如:

/usr/lib/nvidia/sriov-manage -e 0:81:0.0 [root@kvm ~]# /usr/lib/nvidia/sriov-manage -e 0:81:0.0 Enabling VFs on 0000:81:00.0 [root@kvm ~]# Kernel log 可以看到:

		.00.0 has soleware scheduler ENABLED with policy volo_heleAlive.
nvidia	0000:81:00.4:	enabling device (0000 -> 0002)
nvidia	0000:81:00.4:	MDEV: Registered
nvidia	0000:81:00.5:	enabling device (0000 -> 0002)
nvidia	0000:81:00.5:	MDEV: Registered
nvidia	0000:81:00.6:	enabling device (0000 -> 0002)
nvidia	0000:81:00.6:	MDEV: Registered
nvidia	0000:81:00.7:	enabling device (0000 -> 0002)
nvidia	0000:81:00.7:	MDEV: Registered
nvidia	0000:81:01.0:	enabling device (0000 -> 0002)
nvidia	0000:81:01.0:	MDEV: Registered
nvidia	0000:81:01.1:	enabling device (0000 -> 0002)
nvidia	0000:81:01.1:	MDEV: Registered
nvidia	0000:81:01.2:	enabling device (0000 -> 0002)
nvidia	0000.81.01 2.	MDEV. Registered

如果执行 SRIOV 使能时发生如下错误,需要检查 kernel 参数中的 iommu 选项是否生效。 或服务器 BIOS 中的相关选项。

```
[root@kvm nvidia]# cd /usr/lib/nvidia/
[root@kvm nvidia]# ls
common.sh post-install pre-uninstall sriov-manage systemd sysv upstart
[root@kvm nvidia]# lspci |grep -i NV
81:00.0 3D controller: NVIDIA Corporation Device 20f1 (rev a1)
84:00.0 3D controller: NVIDIA Corporation TU104GL [Tesla T4] (rev a1)
[root@kvm nvidia]# ./sriov-manage -e 0:81:0.0
Enabling VFs on 0000:81:00.0
./sriov-manage: line 97: /sys/bus/pci/drivers/nvidia/bind: No such file or directory
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
./sriov-manage: line 142: echo: write error: No such device
```

看到 dmesg kenrel log 有报错:

220 CZOE071 M	VDM. Aborting probe for VE 0000.01.00 A since TOMMU is not present on the sustan
328.670507 N	VRM: Aborting probe for VF 0000:81:00.4 since lommo is not present on the system.
328.670514] n	vidia: probe of 0000:81:00.4 failed with error -1
328.670668] N	VRM: Aborting probe for VF 0000:81:00.5 since IOMMU is not present on the system.
328.670673] n	vidia: probe of 0000:81:00.5 failed with error -1
328.670807] N	VRM: Aborting probe for VF 0000:81:00.6 since IOMMU is not present on the system.
328.670811] n	vidia: probe of 0000:81:00.6 failed with error -1
328.670941] N	VRM: Aborting probe for VF 0000:81:00.7 since IOMMU is not present on the system.
328.670945] n	vidia: probe of 0000:81:00.7 failed with error -1
328.671094] N	VRM: Aborting probe for VF 0000:81:01.0 since IOMMU is not present on the system.
328.671129] n	vidia: probe of 0000:81:01.0 failed with error -1
328.671306] N	VRM: Aborting probe for VF 0000:81:01.1 since IOMMU is not present on the system.
328.671323] n	vidia: probe of 0000:81:01.1 failed with error -1

确认启用了 SR-IOV 之后,对于 A100 和 A30 等支持 MIG 特性的物理 GPU,既可以使用常规 Time Sliced 类型的 vGPU,也可以使用 MIG 后端类型 vGPU。而 A10, A40 等不支持 MIG 特性的 GPU 则只支持 Time Sliced 方式的 vGPU.

通过 lspci 命令可查看 GPU 支持的 SR-IOV 属性,例如查看本案中 81:00.0 的 A100 GPU:

lspci -s 00:81:00.0 -vv



看到 A100 最大分配 VF 数量为 16 个。

所以, 使能 SRIOV 后, 在 /sys/class/mdev\_bus/ 目录中可以列出所有可用于创建 vGPU 的 VF 设备的 BDF, A100 的 00:81:\* 的 VF 目录一共是 **16** 个。

[root@kvm mdev_bus	s]# pwd				
/sys/class/mdev_bu	IS				
[root@kvm mdev_bus	s]# ls				
0000:81:00.4 0000	8:81:00.7	0000:81:01.2	0000:81:01.5	0000:81:02.0	0000:81:02.3
0000:81:00.5 0000	0:81:01.0	0000:81:01.3	0000:81:01.6	0000:81:02.1	0000:84:00.0
0000:81:00.6 0000	9:81:01.1	0000:81:01.4	0000:81:01.7	0000:81:02.2	
[root@kvm mdev_bus	;]#				

#### 2.5. 基于 SRIOV 的 vGPU 特性

创建 SR-IOV 类型的 vGPU, 遵循以下方法:

- 1. 每个 vGPU 实例占用一个 VF 设备,一旦 VF 已经被分配,该 VF 上不可再创建 vGPU。
- 每物理 GPU 的 VF 总数 >= 该 GPU 上可创建的 vGPU 最大实例数,例如 A100 单卡的 Total VF 数量为 16,而最大 A100 单卡的 vGPU 实例数为 40GB(FB 总量)/4GB(4C 类型 vGPU Size) = 10。因此,VF 总数并不是该 GPU 的可创建最大 vGPU 数量。
- 3. 最大可创建的 vGPU 实例数可以查询该 vGPU 类型目录中的 description 文件中描述的 max\_instances 的值。例如:

```
[root@kvm 0000:81:00.4]# cd mdev_supported_types/
[root@kvm mdev_supported_types]# ls
nvidta-468 nvidia-470 nvidia-472 nvidia-474 nvidia-476 nvidia-478
nvidia-469 nvidia-471 nvidia-473 nvidia-475 nvidia-477
[root@kvm mdev_supported_types]# pwd
/sys/class/mdev_bus/0000:81:00.4/mdev_supported_types
[root@kvm mdev_supported_types]# cd nvidia-468/
[root@kvm nvidia-468]# cat name
GRID_A100-4C
[root@kvm nvidia-468]# cat description
num_heads=1, frl_config=60, framebuffer=4096M, max_resolution=4096x2160, max_instance=10
[root@kvm nvidia-468]#
```

4. 每一个 VF 设备的当前实时可创建 vGPU 数量应查询该 VF 目录下,指定类型目录中 available\_instances 文件中的值: 1 表示可以在此 VF 设备创建指定的 vGPU 类型实 例。0 则表示不可创建此类型 vGPU,可能是由于当前 VF 已经被占用,或当前 vGPU 类型不被支持。



下面将分开介绍创建 MIG 和 Non-MIG 模式的 vGPU 实例

如果您使用的是 A10、A40、A16 等非 MIG GPU, 请直接跳至本文第 4 章。

## 3. MIG 类型的 vGPU 配置概要

如果您的物理 GPU 支持 MIG 特性,例如 A100,A30。则既可以使用 MIG 模式的 vGPU 以实现物理 分区的实例隔离性,也可以使用 Time Sliced(时间片方式切分共享)模式,本章节介绍 MIG 模式的 vGPU 的配置。如果您使用时间片切分的 vGPU 请参考下一个章节。

#### 3.1. 启用 MIG 模式

创建基于 MIG 的 vGPU 实例需要将 GPU 切换到 MIG 模式, 此配置为持久配置, 只需执行一次, 重启服务器后仍然有效:

nvidia-smi -mig 1

[root@kvm ~]# nvidia-smi -i @ Enabled MIG Mode for GPU 000@ All done. [root@kvm ~]# nvidia-smi Thu Nov 12 23:39:59 2020	0 -mig 1 00000:81:00.0			
NVIDIA-SMI 450.89 Dri	iver Version: 4	50.89	CUDA Version	: N/A
GPU Name Persistend   Fan Temp Perf Pwr:Usage/   	ce-M  Bus-Id /Cap  M 	Disp.A emory-Usage	Volatile U   GPU-Util   +===================================	Incorr. ECC   Compute M.   MIG M.
0 A100-PCIE-40GB Or N/A 40C P0 87W / 25	n   00000000: 50W   0MiB	81:00.0 Off / 40537MiB	   N/A	On Default Enabled
1 Tesla T4 Or   N/A 38C P8 16W / 7 	n   000000000 70W   83MiB 	84:00.0 Off / 15359MiB	   0% 	0   Default   N/A
+				+
MIG devices:	+	+		 +
GPU GI CI MIG     ID ID Dev   	Memory-Usage   BAR1-Usage	Vol  SM Unc  ECC	Sha CE ENC DE	ired   C OFA JPG  
No MIG devices found	+	=======+		
•				
Processes:   GPU GI CI PID   ID ID	Type Proces	s name		GPU Memory   Usage
No running processes found	d			

3.2. MIG 设备初始化

SR-IOV GPU 的每一个 VF 设备目录下只能创建一个 vGPU 实例。这个在 A100 时分模式的 SRIOV based vGPU 同样适用。

这里我们可以检索每一个 VF 目录所能创建的 vGPU 类型和数量:

[roc	ot@kvm mdev_bus]	# ls							
0000	:81:00.4 0000:	81:00.7	0000:81:	01.2	0000:81:01.5	0000:81:02.0	0000:81:02.3		
0000	:81:00.5 0000:	81:01.0	0000:81:	01.3	0000:81:01.6	0000:81:02.1	0000:84:00.0		
0000	:81:00.6 0000:	81:01.1	0000:81:	01.4	0000:81:01.7	0000:81:02.2			
[roc	t@kvm mdev_bus]	# cd 000	0\:81\:00	.4/mde	v_supported_t	ypes/			
[roc	t@kvm mdev_supp	orted_ty	pes]# ls						
nvid	lia-468 nvidia-	470 nvi	dia-472	nvidia	-474 nvidia-	476 nvidia-47			
nvid	lia-469 nvidia-	471 nvi	dia-473	nvidia	-475 nvidia-	477			
[roc	ot@kvm mdev supp	orted ty	pes]#_for	i in	* ; do echo "	" \$(cat \$i/n	ame) available:	<pre>\$(cat \$i/ava</pre>	ai*); done
G	GRID A100-4C ava	ilable: 0	9						
G	GRID A100-5C ava	ilable: 0	9						
G	GRID A100-8C ava	ilable: 0	9						
G	GRID A100-10C av	ailable:	0						
G	GRID A100-20C av	ailable:	0						
G	GRID A100-40C av	ailable:	0						
G	GRID A100-1-5C a	vailable	: 0						
G	GRID A100-2-10C	availabl	e: 0						
G	GRID A100-3-20C	availabl	e: 0						
G	GRID A100-4-20C	availabl	e: 0						
G	GRID A100-7-40C	availabl	e: 0						
[roc	ot@kvm mdev_supp	orted_ty	pes]#						

可以看到如果不创建 MIG GI,则所以的 MIG 类型的 vGPU 可创建数量均为 0,因此我们需要先为 vGPU 创建相应的 MIG GPU Instance.

有关 MIG 的管理, 请参见:

https://docs.nvidia.com/datacenter/tesla/mig-user-guide/

https://docs.nvidia.com/datacenter/tesla/pdf/NVIDIA\_MIG\_User\_Guide.pdf

3.3. 创建用于 vGPU 的 MIG GPU Instance (GI)

这里我们创建 2 个测试 GI, 分别为 MIG 4g.20gb 和 MIG 2g.10gb

#### 列出 MIG GI Profile:

[r	root@	kvm ~]# nvi	dia-smi	mig -lgip						
	GPU GPU	instance pr Name	ofiles: ID	Instances Free/Total	Memory GiB	P2P	SM CE	DEC JPEG	ENC OFA	-+
	0	MIG 1g.5gb	19	7/7	4.75	No	14 1	0 0	0 0	-     
	0	MIG 2g.10g	jb 14	3/3	9.75	No	28 2	1 0	0 0	
	0	MIG 3g.20g	;b 9	2/2	19.62	No	42 3	2 0	0 0	
İ	0	MIG 4g.20g	;b 5	1/1	19.62	No	56 4	2 Ø	0 0	
	0	MIG 7g.40g	;b 0	1/1	39.50	No	98 7	5 1	0 1	

#### 成功创建 GI:

[r Su Su [r	oot@ cces cces oot@	kvm ~]# sfully c sfully c kvm ~]#	nvidia- created created nvidia-	Smi GPU GPU Smi	mig insta insta mig	-cgi ance ance -lgi	5,1 ID ID	4 1 5	on on	GPU GPU	0 0	using using	g profile g profile	MIG MIG	4g.20gb 2g.10gb	(ID (ID	5) 14)
	GPU GPU	instance Name	es:	Pro I	ofile D	Ins [	stan [D	ce	F	Place Start	mer :Si	nt   ize					
	0	MIG 2g	.10gb	1	4		5			4:	2						
	0	MIG 4g	.20gb		5		1			0:	4						

## 3.4. 创建 vGPU 设备

此时杳看 vGPU profile 各类型中可创建的数量:

[root@kvm mdev_supported_types]# pwd
/sys/class/mdev_bus/0000:81:00.4/mdev_supported_types
[root@kvm mdev_supported_types]# for i in * ; do echo \$i: \$(cat \$i/name) available: \$(cat \$i/avai*); done
nvidia-468: GRID A100-4C available: 0
nvidia-469: GRID A100-5C available: 0
nvidia-470: GRID A100-8C available: 0
nvidia-4/1: GRID A100-10C available: 0
NVIGIA-4/2: GRID A100-20C AVAILADIE: 0
nvidia-473, GRID A100-40C available, o
$r_1 r_1 r_2 r_4 r_5$ (RID A100-2-10C available) 1
vidia 476: GRID A100-3-20C available: 0
vidia-477: GRID A100-4-20C available: 1
nvidia-478: GRID A100-7-40C available: 0
[root@kvm mdev_supported_types]# ls
nvidia-468 nvidia-470 nvidia-472 <u>nvidia-474</u> nvidia-476 nvidia-478
nvidia-469 nvidia-471 nvidia-473 nvidia-475 nvidia-477
即可对于 A100-2-10C,或者 A100-4-20C 创建 MIG vGPU 实例。
这里选择第一个 VF 目录 /sys/class/mdev_bus/0000:81:00. <mark>4</mark> /mdev_supported_types
分别对应目录 nvidia-475, nvidia-477。
[root@kvm mdev supported types]# cd nvidia-475/
[root@kvm nvidia-475]# ls
available_instances create description device_api devices name
[root@kvm nvidia-475]# uuidgen > create
[root@kvm nvidia-475]# ls
available_instances create description device_api devices name
[root@kvm n∀idia-475]# ls devices
c9035e18-77ac-4b54-9df0-284ff30feece
[root@kvm nvidia-475]# _

上面创建了一个 2-10C vGPU 设备,如果创建第二个,则需要到另一个 VF 目录,这里转到 /sys/class/mdev\_bus/0000:81:00.5/mdev\_supported\_types。可以看到仅有 1 个 4-20C 的 vGPU 类型可以创建。



### 3.5. 创建 VM, 向 VM 添加 MIG vGPU 设备

- 1. 创建 VM,并添加 vGPU mdev 设备
- 2. 此步与 Non-MIG 模式 vGPU 的配置相同, 这里请参见标准添加步骤。 以下是 KVM 虚拟机中添加的 mdev 设备的片段。



#### 3.6. VM 内安装 vGPU 客户机驱动程序和 CUDA

VM 启动后,需要在虚拟机内安装 vGPU Guest 驱动程序和 CUDA 框架。

Installing vGPU gues	st driver, ple	ease wait.	•••			
Uncompressing NVIDI	A Accelerated	Graphics [	Driver for Li	nux-x86_64	450.89	
Fri Nov 13 00:45:18	2020					
NVIDIA-SMI 450.89	Driver	Version: 4	450.89	CUDA Versi	lon: 11.0	+
GPU Name     Fan Temp Perf   	Persistence-M Pwr:Usage/Cap	Bus-Id I	Disp.A Memory-Usage	Volatile   GPU-Util 	e Uncorr. ECC Compute M. MIG M.	
0 GRID A100-4-:   N/A N/A P0 	20C On   N/A / N/A	00000000 1844Mil	:00:0A.0 Off B / 20475MiB	+     N/A 	On Default Enabled	
+				+		+
MIG devices:	·		L			
GPU GI CI MIG   ID ID Dev	Memc   BA	ory-Usage AR1-Usage	Vol  SM Unc  ECC	S CE ENC	Shared DEC OFA JPG	
No MIG devices for	+=====================================		+======+			
+						+
Processes:   GPU GI CI   ID ID	PID Typ	e Proces	ss name		GPU Memory Usage	
No running proces	sses found					
Installing CUDA, plo Done. Building gpu_burn [root@test000081005	ease wait 4-20C 1 ~]#					Ŧ

A100 MIG 模式的 vGPU 需要 vCS license 授权

## 3.7. VM 内创建 MIG Compute Instance

MIG vGPU 需要创建计算实例才能运行 CUDA 计算: 也可以根据需要创建不同规格的 CI 实例,默认会创建最大 profile 的 MIG 计算实例。 nvidia-smi mig -cci root@test000081005\_4-20C\_1 ~]# nvidia-smi mig -cci successfully created compute instance ID 0 on GPU 0 GPU instance ID 0 using profile MIG 4g.20gb (ID 3) root@test000081005\_4-20C\_1 ~]# nvidia-smi ri Nov 13 00:55:50 2020 NVIDIA-SMI 450.89 Driver Version: 450.89 CUDA Version: 11.0 GPU Name Persistence-M| Bus-Id Disp.A Volatile Uncorr. ECC Temp Perf Pwr:Usage/Cap Memory-Usage GPU-Util Compute M. Fan MIG M. \_\_\_\_\_ GRID A100-4-20C 00000000:00:0A.0 Off On 0 N/A N/A PØ 1844MiB / 20475MiB N/A Default Enabled MIG devices: GPU GI CI MIG ID ID Dev Memory-Usage Vol Shared BAR1-Usage CE ENC DEC OFA JPG 1844MiB / 20475MiB | 0MiB / 4096MiB | 56 Processes: PID Type Process name GPU Memory Usage No running processes found

#### 3. CUDA 计算测试

[root@test000081005\_4-20C\_1 ~]# nvidia-smi -L GPU 0: GRID A100-4-20C (UUID: GPU-98420d02-2505-11eb-8d63-b7a74215b053) MIG 4g.20gb Device 0: (UUID: MIG-GPU-98420d02-2505-11eb-8d63-b7a74215b053/0/0) root@test000081005\_4-20C\_1 ~]# ls in nvdrv-install.sh NVIDIA\_CUDA-11.0\_Samples NVIDIA-Linux-x86\_64-450.89-grid.run update\_xpciid.sh [root@test000081005\_4-20C\_1 ~]# cd bin/ root@test000081005\_4-20C\_1 bin]# ls compare.ptx gpu\_burn root@test000081005\_4-20C\_1 bin]# ./gpu\_burn 20 GPU 0: GRID A100-4-20C (UUID: GPU-98420d02-2505-11eb-8d63-b7a74215b053) MIG 4g.20gb Device 0: (UUID: MIG-GPU-98420d02-2505-11eb-8d63-b7a74215b053/0/0) Initialized device 0 with 20475 MB of memory (18183 MB available, using 16364 MB of it), using FLOATS 15.0% proc'd: 1020 (5605 Gflop/s) errors: 0 temps: --Summary at: Fri Nov 13 00:59:20 CST 2020 35.0% proc'd: 3060 (8761 Gflop/s) errors: 0 temps: --Summary at: Fri Nov 13 00:59:24 CST 2020 50.0% proc'd: 4080 (8761 Gflop/s) errors: 0 temps: --Fri Nov 13 00:59:27 CST 2020 Summary at: 65.0% proc'd: 6120 (8761 Gflop/s) errors: 0 temps: --Fri Nov 13 00:59:30 CST 2020 Summary at: proc'd: 8160 (8761 Gflop/s) errors: 0 temps: --85.0% Summary at: Fri Nov 13 00:59:34 CST 2020 100.0% proc'd: 9180 (8761 Gflop/s) errors: 0 temps: --Summary at: Fri Nov 13 00:59:37 CST 2020

看到运行 CUDA 计算正常。

## 3.8. 修改 MIG Compute Instance 并为不同的应用指定不同的 CI

• 列出当前分配的 CI:

nvidia-smi -lci

[root@	test00008:	1005_4-20C_1 ~]# n\	/idia-smi mig	g -lci	
Comp	ute insta	nces:	Profilo	Instance	Placement
	Instance	e	ID	ID	Start:Size
  =====	1D ======				======
0	0	MIG 4g.20gb	3	0	0:4

• 删除不需要的 CI:

nvidia-smi -dci -ci <CI\_ID>

[root@test000081005\_4-20C\_1 ~]# nvidia-smi mig -dci -ci 0 Successfully destroyed compute instance ID 0 from GPU 0 GPU instance ID 0

• 列出支持的 CI 类型

nvidia-smi -lcip

[root@t +	est0000810	05_4-20C_1 ~]# nvi	dia-smi mi	g -lcip				
Compu   GPU 	te instanc GPU Instance ID	e profiles: Name	Profile ID	Instances Free/Total	Exclusive SM	DEC CE	Shared ENC JPEG	OFA
======   0 	0	MIG 1c.4g.20gb	0	4/4	14	2 4	0 0	 Ø
+   0 	0	MIG 2c.4g.20gb	1	2/2	28	2 4	0 0	0
0	0	MIG 4g.20gb	3*	1/1	56	2 4	0 0	0

• 可以一次创建多个新的 CI

nvidia-smi -cci <id1,id2,id3...>

下面分别创建 2 个 1c 和一个 2c 实例, Profile ID 分别为 0, 和 1。

[root@test000	081005_4	-20C_1 ~	~]# nvidia	a-smi	i m	nig	-cci	1,	,0,0									
Successfully	created	compute	instance	ID	0	on	GPU	0	GPU	instance	ID	0	using	profile	MIG	2c.4g.20gb	(ID	1)
Successfully	created	compute	instance	ID	1	on	GPU	0	GPU	instance	ID	0	using	profile	MIG	1c.4g.20gb	(ID	0)
Successfully	created	compute	instance	ID	2	on	GPU	0	GPU	instance	ID	0	using	profile	MIG	1c.4g.20gb	(ID	0)
root@test000	081005_4	-20C_1 ~	~]#															

• 设定 CUDA\_VISIBLE\_DEVICES 环境变量,运行 CUDA 程序。



通过上面测试,可以观察到 MIG 模式下的 vGPU VM 内通过 CI 划分,在应用间可以分配不同的资源,提供良好的 QoS 和灵活性。

## 4. Time Sliced 类型的 vGPU 配置概要

对于 Non-MIG 的 Ampere 架构 GPU, vGPU 仍然使用 Time Sliced(时分共享)模式。新架构启用 了 SR-IOV, vGPU 的管理还是在 vfio-mdev 的基础之上的。SR-IOV 的 vGPU 设备分配和之前非 SR-IOV 的 mdev 设备创建和管理会有所不同。本章节将讲解基于 SR-IOV 的 vGPU (例如 A10, A40, A16)时分切分方式的配置方法。

#### 4.1. 创建 vGPU 设备

这里我们仍然以 A100 GPU 为例,在启用 SR-IOV 之后,可以在 /sys/class/mdev\_bus 目录中 看到所有的 VF 设备目录,如下图。

<pre>[root@kvm mdev_bus]# cd /s</pre>	sys/class/mdev_bu	us/			
[root@kvm mdev_bus]# ls					
0000:81:00.4 0000:81:00.3	7 0000:81:01.2	0000:81:01.5	0000:81:02.0 0000	:81:02.3	
0000:81:00.5 0000:81:01.0	0000:81:01.3	0000:81:01.6	0000:81:02.1 0000	:84:00.0	
0000:81:00.6 0000:81:01.3	1 0000:81:01.4	0000:81:01.7	0000:81:02.2		
[root@kvm mdev_bus]# ls 00	300\:81\:00.4				
ari_enabled	d3cold_allowed	iommu	<pre>mdev_supported_</pre>	types resource	subsystem
broken_parity_status	device	iommu_group	modalias	resource0	subsystem_device
class	dma_mask_bits	irq	msi_bus	resource1	subsystem_vendor
config	driver	local_cpulist	numa_node	resource1_wc	uevent
consistent_dma_mask_bits	driver_override	local_cpus	physfn	resource3	vendor
current_link_speed	enable	<pre>max_link_spee</pre>	d power	resource3_wc	
current_link_width	firmware_node	max_link_widt	h reset	revision	
[root@kvm mdev_bus]# _					
[root@kvm mdev_bus]#					

注意:这里基于 SR-IOV 的 vGPU,每一个 VF 设备只对应一个 vGPU 实例。这是非 SRIOV 的 GPU 的主要区别。每个 VF 设备目录,例如上图中 0:81:00.4 VF 中只能创建一个 mdev 设备。同时 Time sliced 模式的 vGPU,仍然要遵循同物理设备的 vGPU 类型唯一的原则。 创建 vGPU 设备仍然是通过每个 VF 目录 /mdev\_supported\_types/\*/create 来创建。 这里我们创建第一个 A100-4C 类型的 vGPU 实例:

1. 选择未使用的 VF 设备目录

选择第一个设备目录: 0000:81:00.4

进入目录可以看到如下图的目录结构:

[root@kvm mdev_bus]#					
[root@kvm mdev_bus]# cd 0	000\:81\:00.4				
[root@kvm 0000:81:00.4]#	ls				
ari_enabled	d3cold_allowed	iommu	<pre>mdev_supported_types</pre>	resource	subsystem
broken_parity_status	device	iommu_group	modalias	resource0	subsystem_device
class	dma_mask_bits	irq	msi_bus	resource1	subsystem_vendor
config	driver	local_cpulist	numa_node	resource1_wc	uevent
consistent_dma_mask_bits	driver_override	local_cpus	physfn	resource3	vendor
current_link_speed	enable	<pre>max_link_speed</pre>	power	resource3_wc	
current_link_width	firmware_node	max_link_width	reset	revision	
[root@kvm 0000:81:00.4]#	cd mdev_supported	_types/			
[root@kvm mdev_supported_	types]# ls				
nvidia-468 nvidia-470 n	vidia-472 nvidia	-474 nvidia-476	nvidia-478		
nvidia-469 nvidia-471 n	vidia-47 <u>3</u> nvidia	-475 nvidia-477			
[root@kvm mdev_supported_	types]#				

2. 找出 4C 类型对应的目录

下面脚本可以列出所有支持的 mdev 设备名称和 vGPU 设备可用数量

cd mdev\_supported\_types for i in \*; do echo \$i, `cat \$i/name` `cat \$i/ava\*`; done [root@kvm mdev\_supported\_types]# [root@kvm mdev\_supported\_types]# for i in \*; do echo \$i, `cat \$i/name` `cat \$i/ava\*`; done nvidia-468, GRID A100-4C 1 nvidia-469, GRID A100-5C 1 nvidia-470, GRID A100-8C 1 nvidia-471, GRID A100-10C 1 nvidia-472, GRID A100-20C 1 nvidia-473, GRID A100-1-5C 0 nvidia-475, GRID A100-1-5C 0 nvidia-475, GRID A100-2-10C 0 nvidia-476, GRID A100-2-0C 0 nvidia-477, GRID A100-2-0C 0 nvidia-478, GRID A100-7-40C 0 [root@kvm mdev\_supported\_types]#

这里可用看到, A100-4C 类型的目录名为 nvidia-468, 且可用实例数量为 1

所以可以在此 VF 上创建 1 个 4C 类型的 vGPU 实例, 目录使用 nvidia-468

3. 使用 uuid 创建 vGPU 设备

uuidgen > nvidia-468/create

[root@kvm mdev\_supported\_types]# cd nvidia-468/ [root@kvm nvidia-468]# ls available\_instances create description device\_api devices name [root@kvm nvidia-468]# uuidgen > create [root@kvm nvidia-468]# ls devices 59b486f4-4205-438f-a24e-2ac2b5384a75

上图可以看到创建成功

[root@kvm mdev_supported_types]# pwd
/sys/class/mdev_bus/0000:81:00.4/mdev_supported_types
[root@kvm mdev_supported_types]#
[root@kvm mdev_supported_types]# for i in * ; do echo \$i, `cat \$i/name` `cat \$i/ava*` ; done
nvidia-468, GRID A100-4C 0
nvidia-469, GRID A100-5C 0
nvidia-470, GRID A100-8C 0
nvidia-471, GRID A100-10C 0
nvidia-472, GRID A100-20C 0
nvidia-473, GRID A100-40C 0
nvidia-474, GRID A100-1-5C 0
nvidia-475, GRID A100-2-10C 0
nvidia-476, GRID A100-3-20C 0
nvidia-477, GRID A100-4-20C 0
nvidia-478, GRID A100-7-40C 0

上图显示已分配 vGPU 的 VF 0:81:00.4 可用的剩余 vGPU 实例数, 全部为 0。意味着当 前 VF 已经不能再创建新的 vGPU 设备。新的 vGPU 只能创建在不同的 VF 上。 因此需要切换到另一个 VF 目录,例如: 0000:81:00.5 检查可用 vGPU 实例状态:

[root@kvm 0000:81:00.5]# pwd	
/sys/class/mdev_bus/0000:81:00.5	
[root@kvm 0000:81:00.5]# cd mdev_supported_types/	
[root@kvm mdev_supported_types]# for iin * ; do echo \$i, `cat \$i/name` `cat \$i/ava*` ; done	
nvidia-468, GRID A100-4C 1	
nvidia-469, GRID A100-5C 0	
nvidia-470, GRID A100-8C 0	
nvidia-471, GRID A100-10C 0	
nvidia-472, GRID A100-20C 0	
nvidia-473, GRID A100-40C 0	
nvidia-474, GRID A100-1-5C 0	
nvidia-475, GRID A100-2-10C 0	
NVIdia-4/6, GRID A100-3-200 0	
NVIDIA-4//, GRID A100-4-200 0	
可以看到全闲的 VF 工面仍然可以创建工作相同类型的 VGPO 关例。直到所有 VF 的可	
用 vGPU instance 全部为 0。	
凸此, 下图在 VF 0.61.00.5 工成功创建 ] 另二十 VGPO 设备。 以此关准。	
[ <u>r</u> oot@kvm mdev supported types]# pwd	
/svs/class/mdev bus/0000.81.00.5/mdev supported types	
[root@kvm_mdev_supported_types]# Is	
nvidia-468 nvidia-470 nvidia-472 nvidia-474 nvidia-476 nvidi	a-478
nvidia 460 nvidia 471 nvidia 472 nvidia 475 nvidia 477	
[root@kvm mdev_supported_types]# pwd	
/sys/class/mdev bus/0000:81:00.5/mdev supported types	
[root@kym mdey supported types]# uuidgen > nyidia-468/create]	
[noot@kym mdoy supported_types]# ]s nyidia_468/dovices/	

43805aeb-9d0a-4c11-bf66-83173eacb660 [root@kvm mdev\_supported\_types]#

#### 4.2. 添加 vGPU MDEV 设备到 VM

后续步骤和标准的非 SR-IOV 的操作相同。

如下图,将之前生成 vGPU 的 uuid 写入虚拟机的 libvirt xml 文件。

```
[root@kvm mig-4-20c]# ls /sys/bus/mdev/devices/ -1
total 0
lrwxrwxrwx 1 root root 0 Jul 13 20:08 30af52db-81c2-487c-adfe-fb5ced0cf806 -> ../../..
0:03.0/0000:82:00.0/0000:83:08.0/0000:84:00.0/30af52db-81c2-487c-adfe-fb5ced0cf806
lrwxrwxrwx 1 root root 0 Aug 5 14:18 43805aeb-9d0a-4c11-bf66-83173eacb660
lrwxrwxrwx 1 root root 0 Aug 5 13:48 59b486f4-4205-438f-a24e-2ac2b5384a75 -> ../../..
0:02.0/0000:81:00.4/59b486f4-4205-438f-a24e-2ac2b5384a75
lrwxrwxrwx 1 root root 0 Jul 13 20:17 f88e4d52-87f7-4d71-86cf-8dd108617db3 -> ../../..
0:03.0/0000:82:00.0/0000:83:08.0/0000:84:00.0/f88e4d52-87f7-4d71-86cf-8dd108617db3
```



之后启动 VM:

```
a100_000081005_4-20C_1 : ca989eb3-56e3-43f4-8e30-f64b0e93fc74
Net:
192.168.122.191/24 | VNET3 | 52:54:00:E2:36:01
Disk:
Type Device Target Source
file disk vda /var/lib/libvirt/images/a100_linked_1-ca989eb3-56e3-43f4-8e30-f64b0e93fc74.qcow2
```

#### 4.3. 安装 vGPU Guest Driver

[root@a100_00081005_4-20C_1 ~]# ls
bin nv-docker-install.sh NVIDIA_CUDA-11.0_Samples NVIDIA-Linux-x86_64-460.73.01-grid.run perl5
[root@a100_000081005_4-20C_1 ~]# ./NVIDIA-Linux-x86_64-460.73.01-grid.run
Verifying archive integrity OK
Uncompressing NVIDIA Accelerated Graphics Driver for Linux-x86_64 460.73.01
[root@a100_00081005_4-20C_1_~]#

安装完成

root@a100_000081005_4-20C_1 ~]# nvidia-smi nu Aug 5 02:32:03 2021
NVIDIA-SMI 460.73.01 Driver Version: 460.73.01 CUDA Version: 11.2
GPU       Name       Persistence-M       Bus-Id       Disp.A       Volatile Uncorr. ECC         Fan       Temp       Perf       Pwr:Usage/Cap       Memory-Usage       GPU-Util       Compute M.         Image: Compute M.       Image       Image       MIG M.
0 GRID A100-4C On   00000000:00:0A.0 Off   0 0 N/A N/A P0 N/A / N/A   407MiB / 4091MiB   0% Default Disabled
Processes:
GPU     GI     CI     PID     Type     Process name     GPU Memory       ID     ID     Usage
No running processes found
~oot@a100 000081005 4-20C 1 ~]#

nvidia-smi 正常显示信息

```
[root@a100_000081005_4-20C_1 ~]# cd bir
[root@a100_000081005_4-20C_1 bin]# ls
 compare.ptx gpu_burn
[root@a100_000801005_4-20C_1 bin]# ./gpu_burn 5
GPU 0: GRID A100-4C (UUID: GPU-01bd2cb9-f5b5-11eb-87db-9d30a27a9dc2)
Initialized device 0 with 4091 MB of memory (3012 MB available, using 2711 MB of it), using FLOATS 20.0% proc'd: 167 (2446 Gflop/s) errors: 0 temps: --
Summary at: Thu Aug 5 02:33:46 EDT 2021
40.0% proc'd: 835 (17087 Gflop/s) errors: 0 temps: --
                          Thu Aug 5 02:33:47 EDT 2021
         Summary at:
60.0% proc'd: 1670 (16978 Gflop/s) errors: 0 temps: --
Summary at: Thu Aug 5 02:33:48 EDT 2021
80.0% proc'd: 2672 (16986 Gflop/s) errors: 0 temps: --
         Summary at: Thu Aug 5 02:33:49 EDT 2021
100.0% proc'd: 3674 (16985 Gflop/s) errors: 0 temps: --
         Summary at: Thu Aug 5 02:33:50 EDT 2021
100.0% proc'd: 4676 (16787 Gflop/s) errors: 0 temps: --
Killing processes.. Freed memory for dev 0
Uninitted cublas
done
Tested 1 GPUs:
         GPU 0: OK
[root@a100 000081005 4-20C 1 bin]#
```

CUDA 测试正常。

4.4. 总结

基于 SR-IOV 的 vGPU 和之前非 SR-IOV 的配置区别主要体现在: 1. 对 SR-IOV 硬件和 Kernel 参数的使能。2. 必须使能 GPU 的 SR-IOV 配置。3. 一个 vGPU 实例对应且仅对应 一个 GPU VF 设备。4. VM 的设备添加、驱动程序安装以及使用与非 SR-IOV 方式相同。 不限于 vCS 类型的 vGPU, 也适用于 vPC 和 vWS。